

Faster Lagrangian-Based Methods in Convex Optimization

Marc Teboulle

School of Mathematical Sciences
Tel Aviv University

Joint work with

Shoham Sabach (Technion)

Optimization without Borders
Dedicated to Yurii Nesterov's 65th Birthday

HSE University, Sochi, Russia, July 12–31, 2021

HAPPY BIRTHDAY Yurii !

Main goal. To unify, simplify, and improve the convergence rate analysis of Lagrangian-based methods for convex optimization.

- A central tool for Lagrangian methods: Nice Primal Algorithmic Map
- A framework of Faster LAGrangian (FLAG) methods
- New non-ergodic rate of convergence results in terms of function values and feasibility violation

We focus on the **linearly constrained** convex optimization problem defined by

$$(P) \quad \min_{x \in \mathbb{R}^n} \{ \Psi(x) : \mathcal{A}x = b \},$$

where

- $\Psi : \mathbb{R}^n \rightarrow (-\infty, +\infty]$ is proper, lsc and σ -strongly convex with $\sigma \geq 0$.
- $\mathcal{A} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is a linear mapping, and $b \in \mathbb{R}^m$.
- The feasible set of problem (P) is denoted by $\mathcal{F} = \{x \in \mathbb{R}^n : \mathcal{A}x = b\} \neq \emptyset$.

Despite its apparent simplicity, this model is very rich and encompasses most convex optimization models.

- **Linear composite model**

$$\min_{u \in \mathbb{R}^p} \{f(u) + g(Au)\} = \min_{u \in \mathbb{R}^p, v \in \mathbb{R}^q} \{f(u) + g(v) : Au = v\},$$

where $f : \mathbb{R}^p \rightarrow (-\infty, +\infty]$ and $g : \mathbb{R}^q \rightarrow (-\infty, +\infty]$ are proper, lower semi-continuous and convex functions, and $A : \mathbb{R}^p \rightarrow \mathbb{R}^q$ is a linear mapping. **It fits into model (P)**, with $x = (u^T, v^T)^T$, $\Psi(x) := f(u) + g(v)$ and $\mathcal{A}x = Au - v$.

- **Linear composite model**

$$\min_{u \in \mathbb{R}^p} \{f(u) + g(Au)\} = \min_{u \in \mathbb{R}^p, v \in \mathbb{R}^q} \{f(u) + g(v) : Au = v\},$$

where $f : \mathbb{R}^p \rightarrow (-\infty, +\infty]$ and $g : \mathbb{R}^q \rightarrow (-\infty, +\infty]$ are proper, lower semi-continuous and convex functions, and $A : \mathbb{R}^p \rightarrow \mathbb{R}^q$ is a linear mapping. **It fits into model (P)**, with $x = (u^T, v^T)^T$, $\Psi(x) := f(u) + g(v)$ and $\mathcal{A}x = Au - v$.

- **Block linear constrained model**

$$\min_{u \in \mathbb{R}^p, v \in \mathbb{R}^q} \{f(u) + g(v) : Au + Bv = b\},$$

where $f : \mathbb{R}^p \rightarrow (-\infty, +\infty]$ and $g : \mathbb{R}^q \rightarrow (-\infty, +\infty]$ are proper, lower semi-continuous and convex functions, $A : \mathbb{R}^p \rightarrow \mathbb{R}^m$ and $B : \mathbb{R}^q \rightarrow \mathbb{R}^m$ are linear mappings. **It fits into model (P)**, with $x = (u^T, v^T)^T$, $\Psi(x) := f(u) + g(v)$ and $\mathcal{A}x = Au + Bv$.

- **Linear composite model**

$$\min_{u \in \mathbb{R}^p} \{f(u) + g(Au)\} = \min_{u \in \mathbb{R}^p, v \in \mathbb{R}^q} \{f(u) + g(v) : Au = v\},$$

where $f : \mathbb{R}^p \rightarrow (-\infty, +\infty]$ and $g : \mathbb{R}^q \rightarrow (-\infty, +\infty]$ are proper, lower semi-continuous and convex functions, and $A : \mathbb{R}^p \rightarrow \mathbb{R}^q$ is a linear mapping. **It fits into model (P)**, with $x = (u^T, v^T)^T$, $\Psi(x) := f(u) + g(v)$ and $\mathcal{A}x = Au - v$.

- **Block linear constrained model**

$$\min_{u \in \mathbb{R}^p, v \in \mathbb{R}^q} \{f(u) + g(v) : Au + Bv = b\},$$

where $f : \mathbb{R}^p \rightarrow (-\infty, +\infty]$ and $g : \mathbb{R}^q \rightarrow (-\infty, +\infty]$ are proper, lower semi-continuous and convex functions, $A : \mathbb{R}^p \rightarrow \mathbb{R}^m$ and $B : \mathbb{R}^q \rightarrow \mathbb{R}^m$ are linear mappings. **It fits into model (P)**, with $x = (u^T, v^T)^T$, $\Psi(x) := f(u) + g(v)$ and $\mathcal{A}x = Au + Bv$.

- **Additive smooth/non-smooth composite objective**

$$\min_{x \in \mathbb{R}^n} \{f(x) + h(x) : \mathcal{A}x = b\},$$

where $f : \mathbb{R}^n \rightarrow (-\infty, +\infty]$ is a proper, lsc and convex function, while $h : \mathbb{R}^n \rightarrow \mathbb{R}$ is a continuously differentiable function with a Lipschitz continuous gradient.

We recall problem (P)

$$(P) \quad \min_{x \in \mathbb{R}^n} \{\Psi(x) : \mathcal{A}x = b\},$$

The corresponding Lagrangian and augmented Lagrangian, are respectively given by:

$$\mathcal{L}(x, y) = \Psi(x) + \langle y, \mathcal{A}x - b \rangle, \quad y \in \mathbb{R}^m,$$

and, for any $\rho > 0$,

$$\mathcal{L}_\rho(x, y) = \mathcal{L}(x, y) + \frac{\rho}{2} \|\mathcal{A}x - b\|^2.$$

Assumption

The Lagrangian \mathcal{L} has a saddle point, that is, there exists (x^, y^*) such that*

$$\mathcal{L}(x^*, y) \leq \mathcal{L}(x^*, y^*) \leq \mathcal{L}(x, y^*), \quad \forall x \in \mathbb{R}^n, \quad \forall y \in \mathbb{R}^m.$$

It can be warranted, for instance, **under standard CQ on the problem's data**.

Starting point: all **Lagrangian-based methods** update a couple (x, y) via

$$\begin{aligned}x^+ &\in \mathcal{P}(x, y), \\y^+ &= y + \mu\rho(\mathcal{A}x^+ - b),\end{aligned}$$

where $\mathcal{P}(\cdot, \cdot)$ is a primal algorithmic map and $\mu > 0$ is a scaling parameter.

Starting point: all **Lagrangian-based methods** update a couple (x, y) via

$$\begin{aligned}x^+ &\in \mathcal{P}(x, y), \\y^+ &= y + \mu\rho(\mathcal{A}x^+ - b),\end{aligned}$$

where $\mathcal{P}(\cdot, \cdot)$ is a primal algorithmic map and $\mu > 0$ is a scaling parameter.

- The main difference between Lagrangian-based methods is encapsulated through the choice of a primal algorithmic map that updates the primal variable.
- This primal map can be seen as the step of any optimization method that is applied on the augmented Lagrangian itself or a variation of it.

$$\mathcal{L}_\rho(x, y) = \Psi(x) + \langle y, \mathcal{A}x - b \rangle + \frac{\rho}{2} \|\mathcal{A}x - b\|^2.$$

Augmented Lagrangian (Hestenes (69), Powell (69))

Main step: Given (x, y) , update the new point (x^+, y^+) via:

$$\begin{aligned}x^+ &\in \operatorname{argmin} \{ \mathcal{L}_\rho(\xi, y) : \xi \in \mathbb{R}^n \}, \\y^+ &= y + \mu\rho (\mathcal{A}x^+ - b).\end{aligned}$$

In this case, \mathcal{P} is an **exact minimization** applied on the augmented Lagrangian.

Examples – Some Classical Lagrangian Based Schemes

$$\mathcal{L}_\rho(x, y) = \Psi(x) + \langle y, \mathcal{A}x - b \rangle + \frac{\rho}{2} \|\mathcal{A}x - b\|^2.$$

Augmented Lagrangian (Hestenes (69), Powell (69))

Main step: Given (x, y) , update the new point (x^+, y^+) via:

$$\begin{aligned}x^+ &\in \operatorname{argmin} \{ \mathcal{L}_\rho(\xi, y) : \xi \in \mathbb{R}^n \}, \\y^+ &= y + \mu\rho (\mathcal{A}x^+ - b).\end{aligned}$$

In this case, \mathcal{P} is an **exact minimization** applied on the augmented Lagrangian.

Proximal Linearized Augmented Lagrangian

Main step: Given (x, y) , update the new point (x^+, y^+) via:

$$\begin{aligned}x^+ &\in \operatorname{argmin} \left\{ \Psi(\xi) + \left\langle \xi, \mathcal{A}^T (y + \rho(\mathcal{A}x - b)) \right\rangle + \frac{1}{2} \|\xi - x\|_M^2 : \xi \in \mathbb{R}^n \right\}, \quad (M \succ 0) \\y^+ &= y + \mu\rho (\mathcal{A}x^+ - b).\end{aligned}$$

In this case, \mathcal{P} is a **proximal gradient** applied on the augmented Lagrangian.

As discussed above, Model (P) covers the following block model

$$\min_{(u,v) \in \mathbb{R}^n} \{f(u) + g(v) : Au + Bv = b\}.$$

$$\mathcal{L}_\rho(u, v, y) = f(u) + g(v) + \langle y, Au + Bv - b \rangle + \frac{\rho}{2} \|Au + Bv - b\|^2.$$

However, the block structure can be exploited in designing Lagrangian-based methods.

As discussed above, Model (P) covers the following block model

$$\min_{(u,v) \in \mathbb{R}^n} \{f(u) + g(v) : Au + Bv = b\}.$$

$$\mathcal{L}_\rho(u, v, y) = f(u) + g(v) + \langle y, Au + Bv - b \rangle + \frac{\rho}{2} \|Au + Bv - b\|^2.$$

However, the block structure can be exploited in designing Lagrangian-based methods.

Alternating Direction Method of Multipliers (ADMM)

(Glowinski and Marroco (75), Gabay and Mercier (76))

Main step: Given (u, v, y) , update the new point (u^+, v^+, y^+) via:

$$u^+ = \operatorname{argmin} \{ \mathcal{L}_\rho(\xi, v, y) : \xi \in \mathbb{R}^n \},$$

$$v^+ = \operatorname{argmin} \{ \mathcal{L}_\rho(u^+, \eta, y) : \eta \in \mathbb{R}^m \},$$

$$y^+ = y + \mu\rho (Au^+ + Bv^+ - b).$$

In this case, \mathcal{P} is an **alternating minimization** applied on the augmented Lagrangian.

$$\mathcal{L}_\rho(u, v, y) = f(u) + g(v) + \langle y, Au + Bv - b \rangle + \frac{\rho}{2} \|Au + Bv - b\|^2.$$

Previous steps can be difficult to implement. Instead, we can *approximate them*:

$$\mathcal{L}_\rho(u, v, y) = f(u) + g(v) + \langle y, Au + Bv - b \rangle + \frac{\rho}{2} \|Au + Bv - b\|^2.$$

Previous steps can be difficult to implement. Instead, we can *approximate them*:

Proximal Linearized ADMM

Main step: Given (u, v, y) , update the new point (u^+, v^+, y^+) via:

$$u^+ = \operatorname{argmin}_\xi \left\{ f(\xi) + \langle A^T (y + \rho(Au + Bv - b)), \xi - u \rangle + \frac{1}{2} \|\xi - u\|_{M_1}^2 \right\},$$
$$v^+ = \operatorname{argmin}_\eta \left\{ g(\eta) + \langle B^T (y + \rho(Au^+ + Bv - b)), \eta - v \rangle + \frac{1}{2} \|\eta - v\|_{M_2}^2 \right\},$$
$$y^+ = y + \mu\rho (Au^+ + Bv^+ - b).$$

(Here $M_1, M_2 \succ 0$).

In this case, \mathcal{P} is a **alternating proximal gradient** applied on the augmented Lagrangian.

Nice Primal Algorithmic Map

It captures the essential ingredients and plays a central role in unifying the analysis of all Lagrangian-based methods into a single and simple framework.

Nice Primal Algorithmic Map

It captures the essential ingredients and plays a central role in unifying the analysis of all Lagrangian-based methods into a single and simple framework.

Definition (Nice primal algorithmic map)

Given the parameters $\rho, t > 0$, let

$$(\rho_t, \tau_t) := \begin{cases} (\rho, t^{-1}) & \text{if } \sigma = 0 \\ (\rho t, t) & \text{if } \sigma > 0. \end{cases}$$

A primal algorithmic map $\text{Prim}_t : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^n$ which is applied on the augmented Lagrangian $\mathcal{L}_{\rho_t}(z, \lambda)$ that generates $z^+ \in \text{Prim}_t(z, \lambda)$, is called **nice**, **if there exist** $\delta \in (0, 1]$ **and** $P, Q \succeq 0$, such that for any $\xi \in \mathcal{F}$ we have

$$\mathcal{L}_{\rho_t}(z^+, \lambda) - \mathcal{L}_{\rho_t}(\xi, \lambda) \leq \tau_t \Delta_P(\xi, z, z^+) - \frac{\tau_t}{2} \|z^+ - z\|_Q^2 - \frac{\sigma}{2} \|\xi - z^+\|^2 - \frac{\delta \rho_t}{2} \|\mathcal{A}z^+ - b\|^2$$

- For any matrix $P \succeq 0$ and any three vectors $u, v, w \in \mathbb{R}^n$:

$$\Delta_P(u, v, w) := \frac{1}{2} \|u - v\|_P^2 - \frac{1}{2} \|u - w\|_P^2.$$

- When $P \equiv I_n$, the identity matrix, we simply write $\Delta_P(u, v, w) \equiv \Delta(u, v, w)$.

FLAG – Faster LAGrangian based method

1. **Input: Problem data** $[\Psi, \mathcal{A}, b, \sigma]$, and a **nice primal algorithmic map** $\text{Prim}_t(\cdot)$.
2. **Initialization:** Set $t_0 = 1$, $\mu \in (0, \delta]$ and $\rho > 0$. Start with any (x^0, z^0, y^0) .

FLAG – Faster LAGrangian based method

1. **Input: Problem data** $[\Psi, \mathcal{A}, b, \sigma]$, and a **nice primal algorithmic map** $\text{Prim}_t(\cdot)$.
2. **Initialization:** Set $t_0 = 1$, $\mu \in (0, \delta]$ and $\rho > 0$. Start with any (x^0, z^0, y^0) .
3. **Iterations:** Generate $\{(x^k, z^k, y^k)\}_{k \in \mathbb{N}}$ and $\{t_k\}_{k \in \mathbb{N}}$ via

3.1. Compute

$$\lambda^k = y^k + \rho_k (t_k - 1) (\mathcal{A}x^k - b), \quad \text{with } \rho_k = \begin{cases} \rho, & \text{if } \sigma = 0 \\ \rho t_k & \text{if } \sigma > 0. \end{cases}$$

3.2. Update the sequence $\{(x^k, z^k, y^k)\}_{k \in \mathbb{N}}$ by

$$z^{k+1} \in \text{Prim}_k(z^k, \lambda^k),$$

$$y^{k+1} = y^k + \mu \rho_k (\mathcal{A}z^{k+1} - b),$$

$$x^{k+1} = (1 - t_k^{-1}) x^k + t_k^{-1} z^{k+1}. \quad \text{The acceleration step!}$$

3.3. Update the sequence $\{t_k\}_{k \in \mathbb{N}}$ by solving the equation $t_{k+1}^p - t_k^p = t_{k+1}^{p-1}$, i.e.,

$$t_{k+1} = \begin{cases} t_k + 1, & p = 1 \text{ (convex case)}, \\ \left(1 + \sqrt{1 + 4t_k^2}\right) / 2, & p = 2 \text{ (strongly convex case)}. \end{cases}$$

- Setting $t_k \equiv 1$ in FLAG, implies $\rho_k \equiv \rho$, $\lambda^k \equiv y^k$, and $x^k \equiv z^k$, thus **recovering the classical basic Lagrangian-based methods**.

- Setting $t_k \equiv 1$ in FLAG, implies $\rho_k \equiv \rho$, $\lambda^k \equiv y^k$, and $x^k \equiv z^k$, thus **recovering the classical basic Lagrangian-based methods**.
- Borrows ideas from the fundamental work on acceleration of Nesterov (1983).
- The choice of t_k plays a key role in accelerating the nice primal map Prim_t . Both the augmented parameter ρ_k and the prox parameter τ_k are determined and chosen through the recursion which defines the sequence t_k .

- Setting $t_k \equiv 1$ in FLAG, implies $\rho_k \equiv \rho$, $\lambda^k \equiv y^k$, and $x^k \equiv z^k$, thus **recovering the classical basic Lagrangian-based methods**.
- Borrows ideas from the fundamental work on acceleration of Nesterov (1983).
- The choice of t_k plays a key role in accelerating the nice primal map Prim_t . Both the augmented parameter ρ_k and the prox parameter τ_k are determined and chosen through the recursion which defines the sequence t_k .
- A main new feature of FLAG is the **auxiliary variable** λ^k defined by:

$$\lambda^k = y^k + \rho_k (t_k - 1) (\mathcal{A}x^k - b),$$

which enable us to derive the new faster **non-ergodic** rate of convergence results!

- Setting $t_k \equiv 1$ in FLAG, implies $\rho_k \equiv \rho$, $\lambda^k \equiv y^k$, and $x^k \equiv z^k$, thus **recovering the classical basic Lagrangian-based methods**.
- Borrows ideas from the fundamental work on acceleration of Nesterov (1983).
- The choice of t_k plays a key role in accelerating the nice primal map Prim_t . Both the augmented parameter ρ_k and the prox parameter τ_k are determined and chosen through the recursion which defines the sequence t_k .
- A main new feature of FLAG is the **auxiliary variable** λ^k defined by:

$$\lambda^k = y^k + \rho_k (t_k - 1) (\mathcal{A}x^k - b),$$

which enable us to derive the new faster **non-ergodic** rate of convergence results!

- Setting $t_k \equiv 1$ in FLAG, implies $\rho_k \equiv \rho$, $\lambda^k \equiv y^k$, and $x^k \equiv z^k$, thus **recovering the classical basic Lagrangian-based methods**.
- Borrows ideas from the fundamental work on acceleration of Nesterov (1983).
- The choice of t_k plays a key role in accelerating the nice primal map Prim_t . Both the augmented parameter ρ_k and the prox parameter τ_k are determined and chosen through the recursion which defines the sequence t_k .
- A main new feature of FLAG is the **auxiliary variable** λ^k defined by:

$$\lambda^k = y^k + \rho_k (t_k - 1) (\mathcal{A}x^k - b),$$

which enable us to derive the new faster **non-ergodic** rate of convergence results!

- Prim_t is assumed to be **nice primal algorithmic map** and this is **all we need** to guarantee rate of convergence results (classical and fast)!

The Two Main Pillars of the Analysis

- The analysis of Lagrangian based methods is usually complicated, and relies on very lengthy and nontrivial proofs.
- Here, it relies on two key lemmas, admitting simple proofs; half-page each!

- The analysis of Lagrangian based methods is usually complicated, and relies on very lengthy and nontrivial proofs.
- Here, it relies on two key lemmas, admitting simple proofs; half-page each!

Lemma 1

Let $\{(x^k, z^k, y^k)\}_{k \in \mathbb{N}}$ generated by FLAG. Then, for any $\xi \in \mathcal{F}$, $\eta \in \mathbb{R}^m$ and $k \geq 0$,

$$\begin{aligned} \mathcal{L}_{\rho_k}(z^{k+1}, \eta) - \mathcal{L}_{\rho_k}(\xi, \eta) &\leq \tau_k \Delta_P(\xi, z^k, z^{k+1}) - \frac{\sigma}{2} \|\xi - z^{k+1}\|^2 + \frac{1}{\mu \rho_k} \Delta(\eta, y^k, y^{k+1}) \\ &\quad - \rho t_{k-1}^p \langle \mathcal{A}x^k - b, \mathcal{A}z^{k+1} - b \rangle. \end{aligned}$$

Proof. Since Prim_k is nice, we have

$$\mathcal{L}_{\rho_k}(z^{k+1}, \lambda^k) - \mathcal{L}_{\rho_k}(\xi, \lambda^k) \leq \tau_k \Delta_P(\xi, z^k, z^{k+1}) - \frac{\sigma}{2} \|\xi - z^{k+1}\|^2 - \frac{\delta \rho_k}{2} \|\mathcal{A}z^{k+1} - b\|^2.$$

From the multiplier update and the three-points identity, we obtain, for all $\eta \in \mathbb{R}^m$,

$$\langle \eta - y^k, \mathcal{A}z^{k+1} - b \rangle = \frac{1}{\mu \rho_k} \langle \eta - y^k, y^{k+1} - y^k \rangle = \frac{1}{\mu \rho_k} \Delta(\eta, y^k, y^{k+1}) + \frac{\mu \rho_k}{2} \|\mathcal{A}z^{k+1} - b\|^2.$$

Using the update rule of the sequence $\{t_k\}_{k \in \mathbb{N}}$ we have that $\rho_k(t_k - 1) = \rho t_{k-1}^p$ and thus

$$\lambda^k = y^k + \rho_k(t_k - 1)(\mathcal{A}x^k - b) = y^k + \rho t_{k-1}^p (\mathcal{A}x^k - b),$$

$$\text{therefore } \langle \eta - \lambda^k, \mathcal{A}z^{k+1} - b \rangle = \langle \eta - y^k, \mathcal{A}z^{k+1} - b \rangle - \rho t_{k-1}^p \langle \mathcal{A}x^k - b, \mathcal{A}z^{k+1} - b \rangle.$$

Combining these relations yields (recall that $\mu \leq \delta$) the desired result. \square

Lemma 2

Let $\{(x^k, z^k, y^k)\}_{k \in \mathbb{N}}$ generated by FLAG. Then, for any $\xi \in \mathcal{F}$, $\eta \in \mathbb{R}^m$ and $k \geq 0$,

$$t_k^p \mathbf{s}_{k+1} - t_{k-1}^p \mathbf{s}_k \leq \frac{\tau_k \rho_k}{\rho} \Delta_P(\xi, z^k, z^{k+1}) - \frac{\rho_k \sigma}{2\rho} \|\xi - z^{k+1}\|^2 + \frac{1}{\mu \rho} \Delta(\eta, y^k, y^{k+1}),$$

where $\mathbf{s}_k = \mathcal{L}_{\rho t_{k-1}^p}(x^k, \eta) - \mathcal{L}_{\rho t_{k-1}^p}(\xi, \eta)$.

Proof of Lemma 2

Proof. Since $x^{k+1} = (1 - t_k^{-1})x^k + t_k^{-1}z^{k+1}$, we obtain from the convexity of $\Psi(\cdot)$ that

$$\Psi(x^{k+1}) \leq (1 - t_k^{-1})\Psi(x^k) + t_k^{-1}\Psi(z^{k+1}).$$

Therefore, multiplying both sides by t_k^p (recalling that $t_k^p - t_k^{p-1} = t_{k-1}^p$), we obtain

$$t_k^p (\Psi(x^{k+1}) - \Psi(\xi)) - t_{k-1}^p (\Psi(x^k) - \Psi(\xi)) \leq t_k^{p-1} (\Psi(z^{k+1}) - \Psi(\xi)).$$

We also have that

$$t_k^p \langle \eta, \mathcal{A}x^{k+1} - b \rangle - t_{k-1}^p \langle \eta, \mathcal{A}x^k - b \rangle = t_k^{p-1} \langle \eta, \mathcal{A}z^{k+1} - b \rangle.$$

Combining these two facts yields

$$t_k^p (\mathcal{L}(x^{k+1}, \eta) - \mathcal{L}(\xi, \eta)) - t_{k-1}^p (\mathcal{L}(x^k, \eta) - \mathcal{L}(\xi, \eta)) \leq t_k^{p-1} (\mathcal{L}(z^{k+1}, \eta) - \mathcal{L}(\xi, \eta)).$$

Now, we will again use that $x^{k+1} = (1 - t_k^{-1})x^k + t_k^{-1}z^{k+1}$ to obtain

$$\|\mathcal{A}x^{k+1} - b\|^2 = (1 - t_k^{-1})^2 \|\mathcal{A}x^k - b\|^2 + t_k^{-2} \|\mathcal{A}z^{k+1} - b\|^2 + 2(1 - t_k^{-1})t_k^{-1} \langle \mathcal{A}x^k - b, \mathcal{A}z^{k+1} - b \rangle.$$

Multiplying both sides of the above equality by $\rho t_k^{2p}/2$ yields

$$\frac{\rho t_k^{2p}}{2} \|\mathcal{A}x^{k+1} - b\|^2 - \frac{\rho t_{k-1}^{2p}}{2} \|\mathcal{A}x^k - b\|^2 = \frac{\rho_k t_k^{p-1}}{2} \|\mathcal{A}z^{k+1} - b\|^2 + \rho_k t_{k-1}^p \langle \mathcal{A}x^k - b, \mathcal{A}z^{k+1} - b \rangle.$$

Therefore, with $s_k = \mathcal{L}_{\rho t_{k-1}^p}(x^k, \eta) - \mathcal{L}_{\rho t_{k-1}^p}(\xi, \eta)$, we have

$$t_k^p s_{k+1} - t_{k-1}^p s_k \leq t_k^{p-1} (\mathcal{L}_{\rho_k}(z^{k+1}, \eta) - \mathcal{L}_{\rho_k}(\xi, \eta)) + \rho_k t_{k-1}^p \langle \mathcal{A}x^k - b, \mathcal{A}z^{k+1} - b \rangle.$$

From the previous lemma, after we multiplied both sides by t_k^{p-1} (recall that $\rho_k = \rho t_k^{p-1}$), we have

$$t_k^{p-1} \left(\mathcal{L}_{\rho_k} \left(z^{k+1}, \eta \right) - \mathcal{L}_{\rho_k} \left(\xi, \eta \right) \right) + \rho_k t_{k-1}^p \left\langle \mathcal{A}x^k - b, \mathcal{A}z^{k+1} - b \right\rangle \leq \\ \frac{\tau_k \rho_k}{\rho} \Delta_P \left(\xi, z^k, z^{k+1} \right) - \frac{\rho_k \sigma}{2\rho} \left\| \xi - z^{k+1} \right\|^2 + \frac{1}{\mu \rho} \Delta \left(\eta, y^k, y^{k+1} \right).$$

By combining the last two inequalities, we obtain the desired result. □

We focus on non-asymptotic rate of convergence (iteration complexity) using the following two classical measures:

- (i) **Function values gap** in terms of $\Psi(x^k) - \Psi(x^*)$.
- (ii) **Feasibility violation** of the constraints of problem (P) in terms of $\|Ax^k - b\|$.

Other measures in the literature: Lagrangian, $\|x^k - x^*\|^2$, $\|x^{k+1} - x^k\|^2$, etc.

When discussing these measures, we also distinguish between rates expressed in terms of the **original produced sequence** or of the **ergodic sequence**.

We focus on non-asymptotic rate of convergence (iteration complexity) using the following two classical measures:

- (i) **Function values gap** in terms of $\Psi(x^k) - \Psi(x^*)$.
- (ii) **Feasibility violation** of the constraints of problem (P) in terms of $\|Ax^k - b\|$.

Other measures in the literature: Lagrangian, $\|x^k - x^*\|^2$, $\|x^{k+1} - x^k\|^2$, etc.

When discussing these measures, we also distinguish between rates expressed in terms of the **original produced sequence** or of the **ergodic sequence**.

- Many rate of convergence results in the literature for variants of Lagrangian-based methods. **Mostly ergodic!** (Chambolle and Pock (11), He and Yuan (12), Monteiro-Svaiter (13),...)

We focus on non-asymptotic rate of convergence (iteration complexity) using the following two classical measures:

- (i) **Function values gap** in terms of $\Psi(x^k) - \Psi(x^*)$.
- (ii) **Feasibility violation** of the constraints of problem (P) in terms of $\|Ax^k - b\|$.

Other measures in the literature: Lagrangian, $\|x^k - x^*\|^2$, $\|x^{k+1} - x^k\|^2$, etc.

When discussing these measures, we also distinguish between rates expressed in terms of the **original produced sequence** or of the **ergodic sequence**.

- Many rate of convergence results in the literature for variants of Lagrangian-based methods. **Mostly ergodic!** (Chambolle and Pock (11), He and Yuan (12), Monteiro-Svaiter (13),...)
- Non-ergodic $O(1/N)$ result for $\|x^{k+1} - x^k\|^2$ (He and Yuan (15)).
- Non-ergodic $O(1/N^2)$ result for $\|x^k - x^*\|^2$, in the strongly convex setting (Chambolle and Pock (11)).

Types of Rate of Convergence – Many Results

We focus on non-asymptotic rate of convergence (iteration complexity) using the following two classical measures:

- (i) **Function values gap** in terms of $\Psi(x^k) - \Psi(x^*)$.
- (ii) **Feasibility violation** of the constraints of problem (P) in terms of $\|Ax^k - b\|$.

Other measures in the literature: Lagrangian, $\|x^k - x^*\|^2$, $\|x^{k+1} - x^k\|^2$, etc.

When discussing these measures, we also distinguish between rates expressed in terms of the **original produced sequence** or of the **ergodic sequence**.

- Many rate of convergence results in the literature for variants of Lagrangian-based methods. **Mostly ergodic!** (Chambolle and Pock (11), He and Yuan (12), Monteiro-Svaiter (13),...)
- Non-ergodic $O(1/N)$ result for $\|x^{k+1} - x^k\|^2$ (He and Yuan (15)).
- Non-ergodic $O(1/N^2)$ result for $\|x^k - x^*\|^2$, in the strongly convex setting (Chambolle and Pock (11)).
- Non-ergodic $O(1/N)$ rate of convergence result in terms of function values and feasibility violation for the **specific Linearized ADMM** (Li and Lin (19)).

The strongly convex case $\sigma > 0$.**Theorem 1. (A fast non-ergodic function values and feasibility violation rates)**

Let $\{(x^k, z^k, y^k)\}_{k \in \mathbb{N}}$ be a sequence generated by FLAG. **Suppose that $\sigma > 0$ and $0 \preceq P \preceq (\sigma/2) I_n$.** Then, for any optimal solution x^* of problem (P) and $N \geq 1$,

$$\Psi(x^N) - \Psi(x^*) \leq \frac{B_{\rho,c}(x^*)}{2N^2} \quad \text{and} \quad \|\mathcal{A}x^N - b\| \leq \frac{B_{\rho,c}(x^*)}{cN^2},$$

where $B_{\rho,c}(x^*) := 4 \left(\|x^* - z^0\|_P^2 + \frac{1}{\mu\rho} (\|y^0\| + c)^2 \right)$.

The convex case $\sigma = 0$.

Theorem 2. (A non-ergodic function values and feasibility violation rates)

Let $\{(x^k, z^k, y^k)\}_{k \in \mathbb{N}}$ be a sequence generated by FLAG and **suppose that $\sigma = 0$** . Then, for any optimal solution x^* of problem (P) and $N \geq 1$,

$$\Psi(x^N) - \Psi(x^*) \leq \frac{B_{\rho,c}(x^*)}{2N} \quad \text{and} \quad \|\mathcal{A}x^N - b\| \leq \frac{B_{\rho,c}(x^*)}{cN},$$

where $B_{\rho,c}(x^*) := 2 \left(\|x^* - z^0\|_P^2 + \frac{1}{\mu\rho} (\|y^0\| + c)^2 \right)$.

FLAG is versatile

- Deriving weaker results of ergodic type was not our primary goal.
- Nevertheless, our main framework FLAG easily adapt to that task.
(The sequences $\{x^k\}_{k \in \mathbb{N}}$ and $\{\lambda^k\}_{k \in \mathbb{N}}$ are not used for that scenario!)

See details in the paper.

Corollary 1. (A fast ergodic function values and feasibility violation rates)

Let $\{(z^k, y^k)\}_{k \in \mathbb{N}}$ be a sequence generated by FLAG. **Suppose that $\sigma > 0$ and $0 \preceq P \preceq (\sigma/2) I_n$.** Then, for any optimal solution z^* of problem (P), the following holds for the **ergodic sequence** $\bar{z}^N = t_{N-1}^{-2} \sum_{k=0}^{N-1} t_k z^{k+1}$

$$\Psi(\bar{z}^N) - \Psi(z^*) \leq \frac{B_{\rho,c}(z^*)}{2N^2} \quad \text{and} \quad \|\mathcal{A}\bar{z}^N - b\| \leq \frac{B_{\rho,c}(z^*)}{cN^2},$$

where $B_{\rho,c}(z^*) := 4 \left(\|z^* - z^0\|_P^2 + \frac{1}{\mu\rho} (\|y^0\| + c)^2 \right)$.

Corollary 2. (An ergodic function values and feasibility violation rates)

Let $\{(z^k, y^k)\}_{k \in \mathbb{N}}$ be a sequence generated by FLAG with $\sigma = 0$ and $t_k = 1$ for all $k \in \mathbb{N}$. Then, for any optimal solution z^* of problem (P), the following holds for **the ergodic sequence** $\bar{z}^N = N^{-1} \sum_{k=0}^{N-1} z^{k+1}$

$$\Psi(\bar{z}^N) - \Psi(z^*) \leq \frac{B_{\rho,c}(z^*)}{2N} \quad \text{and} \quad \|\mathcal{A}\bar{z}^N - b\| \leq \frac{B_{\rho,c}(z^*)}{cN},$$

where $B_{\rho,c}(z^*) := 2 \left(\|z^* - x^0\|_P^2 + \frac{1}{\mu\rho} (\|y^0\| + c)^2 \right)$.

The notion of nice algorithmic map is flexible and **easily adapt to the block setting**:

$$\min_{x:=(u,v) \in \mathbb{R}^p \times \mathbb{R}^q = \mathbb{R}^n} \{\Psi(x) := f(u) + g(v) : \mathcal{A}x := Au + Bv = b\}.$$

In the block model, we only need to assume that **either f or g is strongly convex**.

The notion of nice algorithmic map is flexible and **easily adapt to the block setting**:

$$\min_{x:=(u,v) \in \mathbb{R}^p \times \mathbb{R}^q = \mathbb{R}^n} \{\Psi(x) := f(u) + g(v) : \mathcal{A}x := Au + Bv = b\}.$$

In the block model, we only need to assume that **either f or g is strongly convex**.

Definition (Nice primal algorithmic map - Block version)

Given the parameters $\rho, t > 0$, we let $(\rho_t, \tau_t) = (\rho, t^{-1})$ (when $\sigma = 0$) and $(\rho_t, \tau_t) = (\rho t, t)$ (when $\sigma > 0$). A primal algorithmic map $\text{Prim}_t : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^n$, which is applied on the augmented Lagrangian $\mathcal{L}_{\rho_t}(z, \lambda)$, that generates $z^+ = (u^+, v^+)$ via $z^+ \in \text{Prim}_t(z, \lambda)$, is *nice* if there exist $\delta \in (0, 1]$ and $P_1, Q_1 \in \mathbb{S}_+^p$ and $P_2, Q_2 \in \mathbb{S}_+^q$ with $P = (P_1, P_2)$ and $Q = (Q_1, Q_2)$, s.t. for any $(\xi_1, \xi_2) \in \mathcal{F}$

$$\begin{aligned} \mathcal{L}_{\rho_t}(z^+, \lambda) - \mathcal{L}_{\rho_t}(\xi, \lambda) &\leq \frac{1}{t} \Delta_{P_1}(\xi_1, u, u^+) - \frac{1}{2t} \|u^+ - u\|_{Q_1}^2 + \tau_t \Delta_{P_2}(\xi_2, v, v^+) \\ &\quad - \frac{\tau_t}{2} \|v^+ - v\|_{Q_2}^2 - \frac{\sigma}{2} \|\xi_2 - v^+\|^2 - \frac{\delta \rho_t}{2} \|\mathcal{A}z^+ - b\|^2. \end{aligned}$$

• **Note:** Here we use g strongly convex. Hence, only the block v uses τ_t (to cover both convex/strongly convex cases). For the other block u (with only convexity. i.e., $\sigma = 0$), we fixed $\tau_t = t^{-1}$.

- Augmented Lagrangian Methods (classical, proximal, and prox-linearized)
- Alternating Direction Method of Multipliers ADMM
- Proximal ADMM
- Proximal Linearized ADMM
- Chambolle-Pock Method
- Proximal Jacobi Direction Method of Multipliers
- Predictor Corrector Proximal Multipliers

For each method an explicit parameter δ and matrices P, Q can be found!
(See details in paper.)

- Augmented Lagrangian Methods (classical, proximal, and prox-linearized)
- Alternating Direction Method of Multipliers ADMM
- Proximal ADMM
- Proximal Linearized ADMM
- Chambolle-Pock Method
- Proximal Jacobi Direction Method of Multipliers
- Predictor Corrector Proximal Multipliers

For each method an explicit parameter δ and matrices P, Q can be found!
(See details in paper.)

Meaning, they all **admit Nice Primal Algorithmic Map!**

Therefore, our nonergodic convergence rate results can be applied.

In addition, nice primal algorithmic maps, can be also be identified for problems with **composite objective...**

We consider the sum composite model: nonsmooth + smooth objective

$$\min_{x \in \mathbb{R}^n} \{f(x) + h(x) : \mathcal{A}x = b\},$$

- $f : \mathbb{R}^n \rightarrow (-\infty, +\infty]$ is a proper, lsc and σ -strongly convex ($\sigma \geq 0$),
- $h : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex C^1 with **L -Lipschitz continuous gradient**.

Model (P) with the Sum Composite Objective Function

We consider the sum composite model: nonsmooth + smooth objective

$$\min_{x \in \mathbb{R}^n} \{f(x) + h(x) : \mathcal{A}x = b\},$$

- $f : \mathbb{R}^n \rightarrow (-\infty, +\infty]$ is a proper, lsc and σ -strongly convex ($\sigma \geq 0$),
- $h : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex C^1 with **L-Lipschitz continuous gradient**.

Lemma (Proximal AL is nice)

Let $M \succeq L I_n$, the primal algorithmic map $\text{Prim}_t(\cdot)$ defined by

$$z^+ = \operatorname{argmin}_{\xi} \left\{ f(\xi) + \langle \nabla h(z), \xi \rangle + \langle \lambda, \mathcal{A}\xi - b \rangle + \frac{\rho t}{2} \|\mathcal{A}\xi - b\|^2 + \frac{\tau t}{2} \|\xi - z\|_M^2 \right\},$$

is nice with $\delta = 1$ and $P = M$ and $Q = M - L I_n$.

Model (P) with the Sum Composite Objective Function

We consider the sum composite model: nonsmooth + smooth objective

$$\min_{x \in \mathbb{R}^n} \{f(x) + h(x) : \mathcal{A}x = b\},$$

- $f : \mathbb{R}^n \rightarrow (-\infty, +\infty]$ is a proper, lsc and σ -strongly convex ($\sigma \geq 0$),
- $h : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex C^1 with **L-Lipschitz continuous gradient**.

Lemma (Proximal AL is nice)

Let $M \succeq L I_n$, the primal algorithmic map $\text{Prim}_t(\cdot)$ defined by

$$z^+ = \operatorname{argmin}_{\xi} \left\{ f(\xi) + \langle \nabla h(z), \xi \rangle + \langle \lambda, \mathcal{A}\xi - b \rangle + \frac{\rho t}{2} \|\mathcal{A}\xi - b\|^2 + \frac{\tau t}{2} \|\xi - z\|_M^2 \right\},$$

is nice with $\delta = 1$ and $P = M$ and $Q = M - L I_n$.

Lemma (Proximal Linearized AL is nice)

Let $M \succeq \rho \mathcal{A}^T \mathcal{A} + L I_n$, the primal algorithmic map $\text{Prim}_t(\cdot)$ defined by

$$z^+ = \operatorname{argmin}_{\xi} \left\{ f(\xi) + \langle \nabla h(z), \xi \rangle + \langle \lambda, \mathcal{A}\xi - b \rangle + \rho t \langle \mathcal{A}z - b, \mathcal{A}\xi \rangle + \frac{\tau t}{2} \|\xi - z\|_M^2 \right\},$$

is nice with $\delta = 1$ and $P = M - \rho \mathcal{A}^T \mathcal{A}$ and $Q = M - \rho \mathcal{A}^T \mathcal{A} - L I_n$.

A Simple Recipe for Rate of Convergence of Lagrangian-based Methods

- (i) Formulate the problem at hand via model (P), *i.e.*, **identify the relevant problem data** $[\Psi, \mathcal{A}, b, \sigma]$. **The value of σ will determine the type of rate that can be achieved (classical or fast).**

- (ii) **Define the desired iterative step(s) of the primal algorithmic map** $\text{Prim}_t(\cdot)$ applied on the augmented Lagrangian $\mathcal{L}_{\rho_t}(\cdot)$ of model (P).

A Simple Recipe for Rate of Convergence of Lagrangian-based Methods

- (i) Formulate the problem at hand via model (P), *i.e.*, **identify the relevant problem data** $[\Psi, \mathcal{A}, b, \sigma]$. **The value of σ will determine the type of rate that can be achieved (classical or fast).**
- (ii) **Define the desired iterative step(s) of the primal algorithmic map** $\text{Prim}_t(\cdot)$ applied on the augmented Lagrangian $\mathcal{L}_{\rho_t}(\cdot)$ of model (P).
- (iii) **Show that the defined primal algorithmic map is nice**, *i.e.*, determine the parameter δ and the matrices P and Q .

A Simple Recipe for Rate of Convergence of Lagrangian-based Methods

- (i) Formulate the problem at hand via model (P), *i.e.*, **identify the relevant problem data** $[\Psi, \mathcal{A}, b, \sigma]$. **The value of σ will determine the type of rate that can be achieved (classical or fast).**
- (ii) **Define the desired iterative step(s) of the primal algorithmic map** $\text{Prim}_t(\cdot)$ applied on the augmented Lagrangian $\mathcal{L}_{\rho_t}(\cdot)$ of model (P).
- (iii) **Show that the defined primal algorithmic map is nice**, *i.e.*, determine the parameter δ and the matrices P and Q .
- (iv) Apply Theorem 1 (if $\sigma > 0$) or Theorem 2 (if $\sigma = 0$) for the corresponding FLAG to **obtain a faster non-ergodic rate of convergence for the designed method.**

Therefore, there is no need any more to enter into the machinery of the proofs!

Sabach, S. and Teboulle, M. Faster Lagrangian-Based Methods in Convex Optimization (November, 2020).

Thank you for listening!

Happy Birthday Yurii !

<http://www.math.tau.ac.il/~teboulle/>